

# Linguistic Analysis, Data Mining, and Clustering to Predict Document Age

Student: Keegan Freeman

Research Professor: Dr. Shishir Shah

## Goal of Research

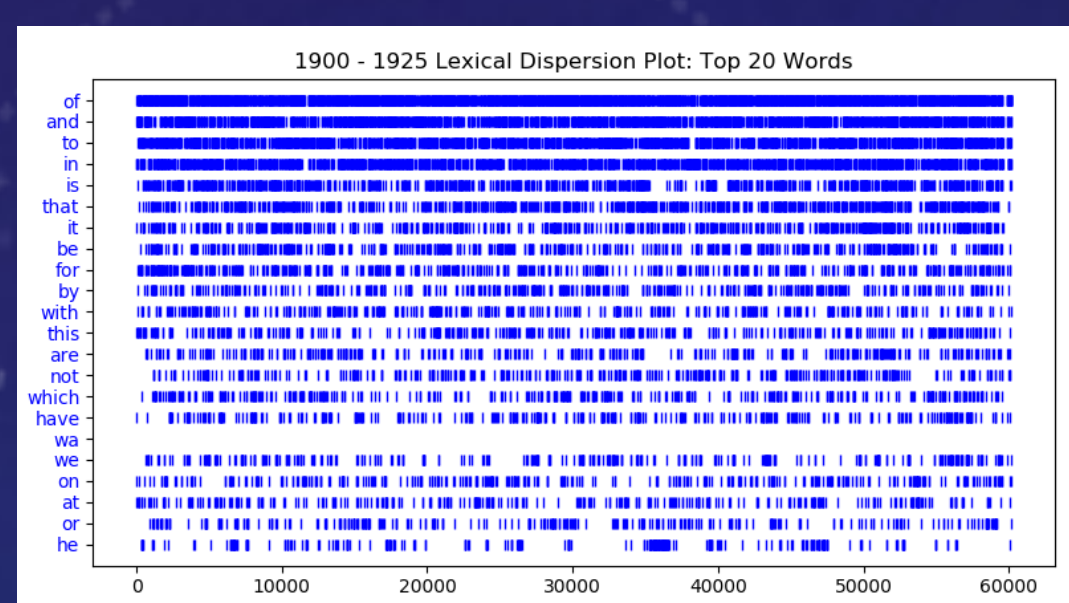
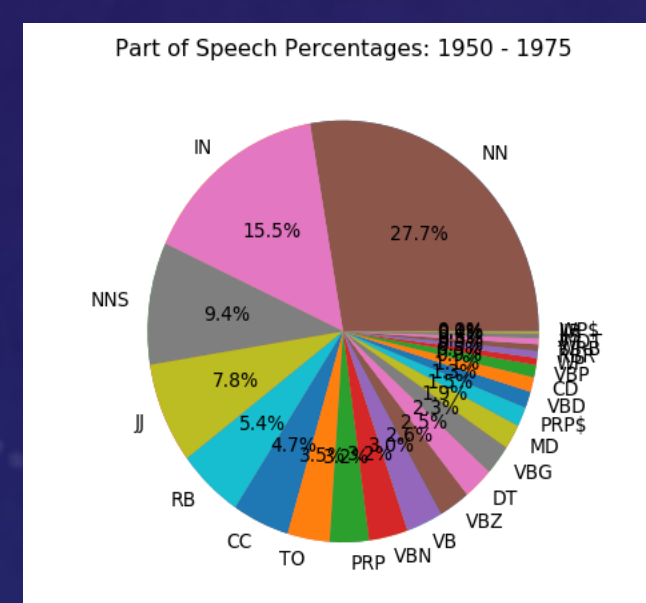
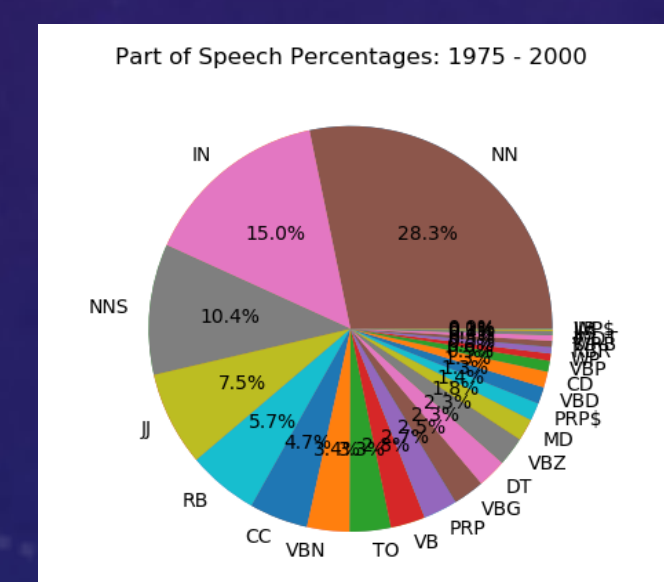
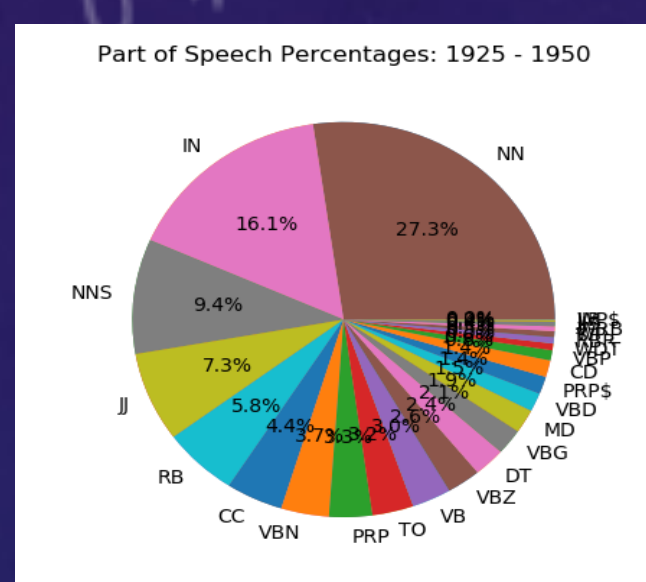
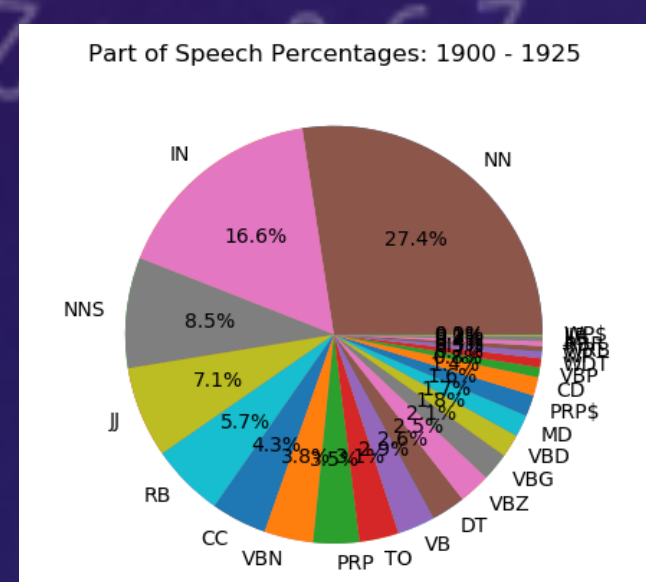
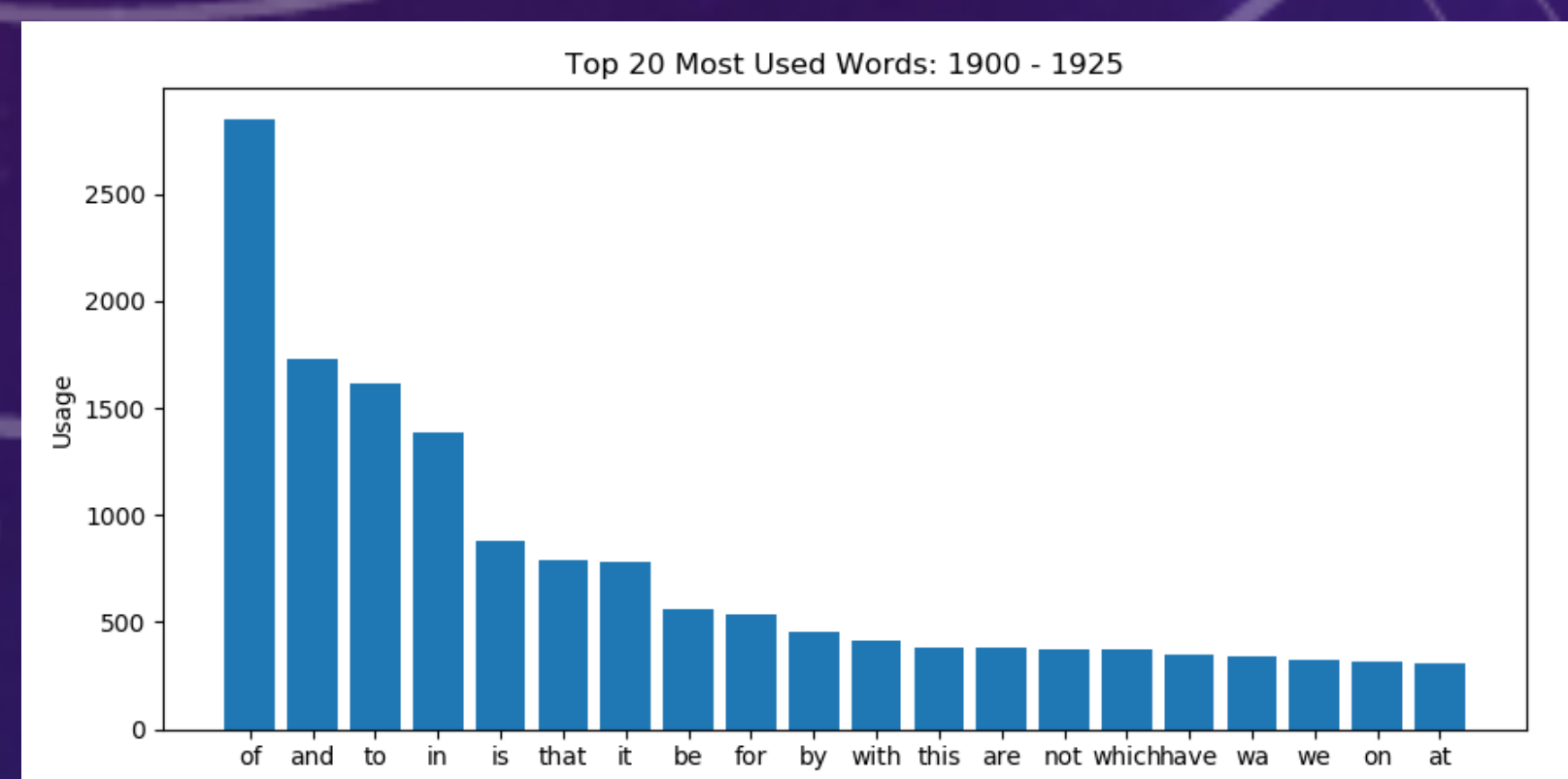
This research aims to create a program capable of iterating through an entire corpora of plain text documents and extract enough linguistic structure data so that, through data clustering, time prediction can be performed on the various documents both in and outside the given corpora. The corpora used is limited to American communication in the form of articles and periodicals less than five pages long so that a diverse corpora representing generalized American communication could be formed.

## Metrics and Collected Data

A variety of data was extracted from each text source and recorded in the written research program. This data includes for each source as well as for each time-span: an array of the stemmed words used, an array of the parts of speech for each stemmed word, word tally, part of speech tally, bigram and trigram collocations for words used together, bigram and trigram collocations for parts of speech used together, word tally percentage, part of speech percentage, and average sentence length.

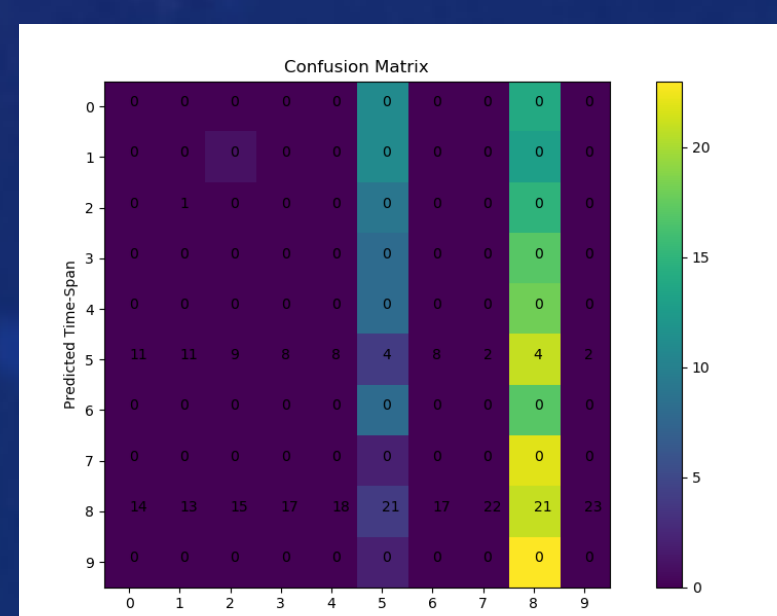
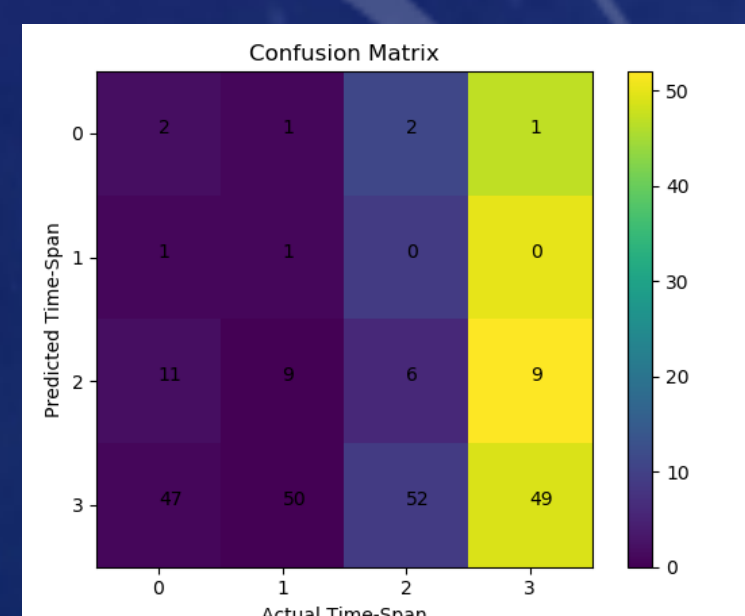
The metrics taken from this collected data to be used in clustering and therefore the data used for time prediction are as follows: average sentence length, top 25 part of speech percentages, top 20 part of speech bigram collocations, and top 20 part of speech trigram collocations.

After storing each text file's metrics into a data point, filtering was performed to deem which metrics are most beneficial to be used for clustering.



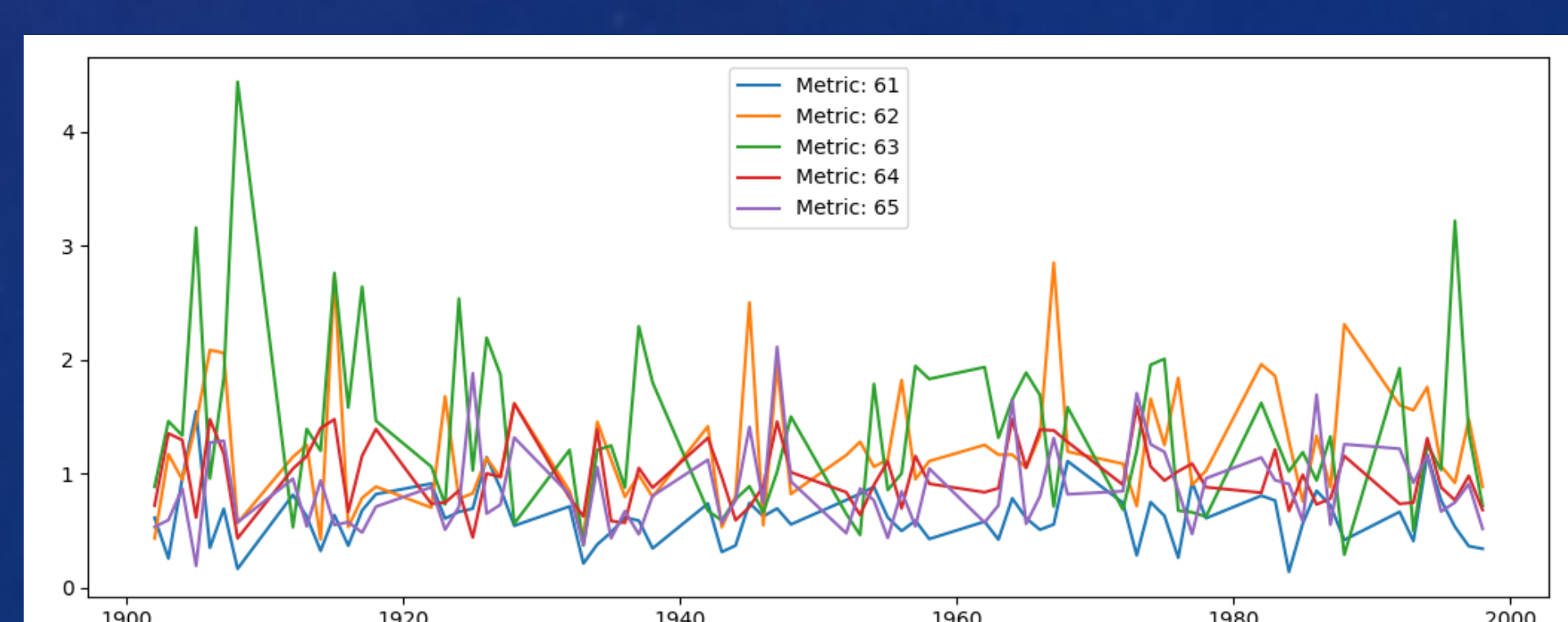
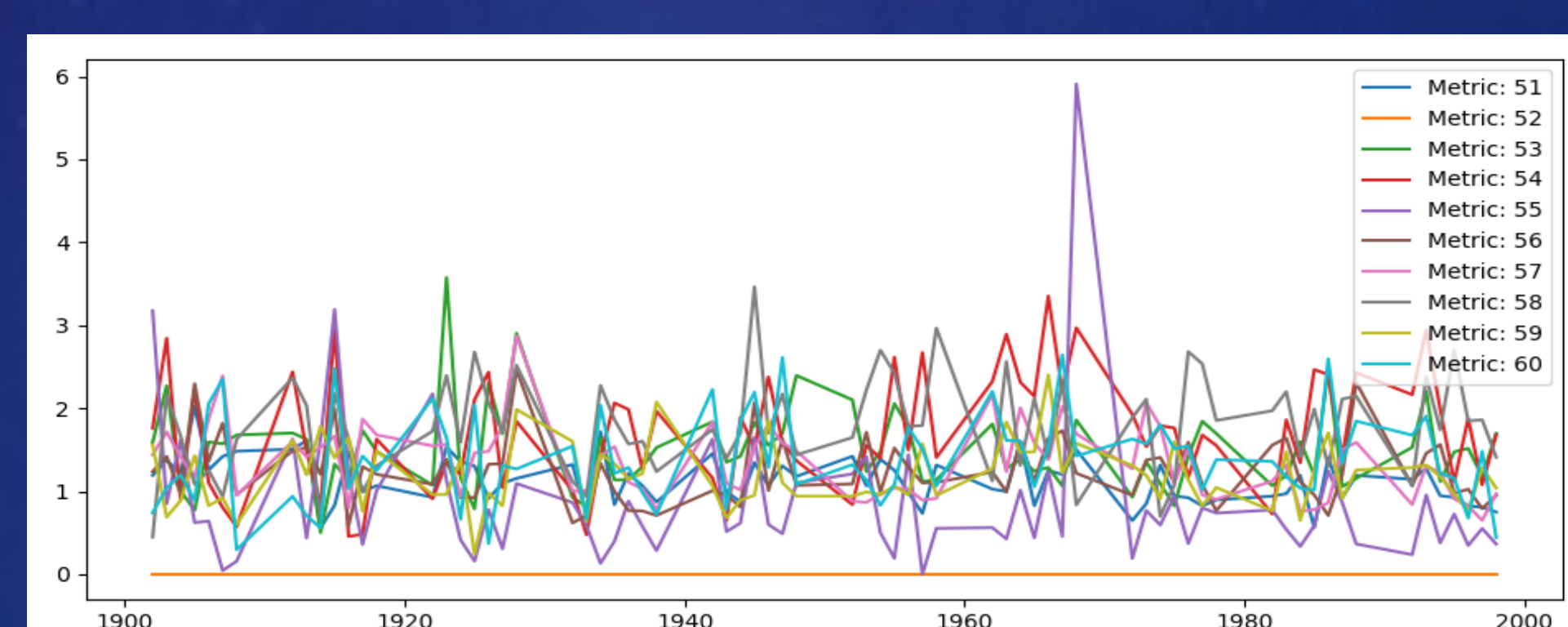
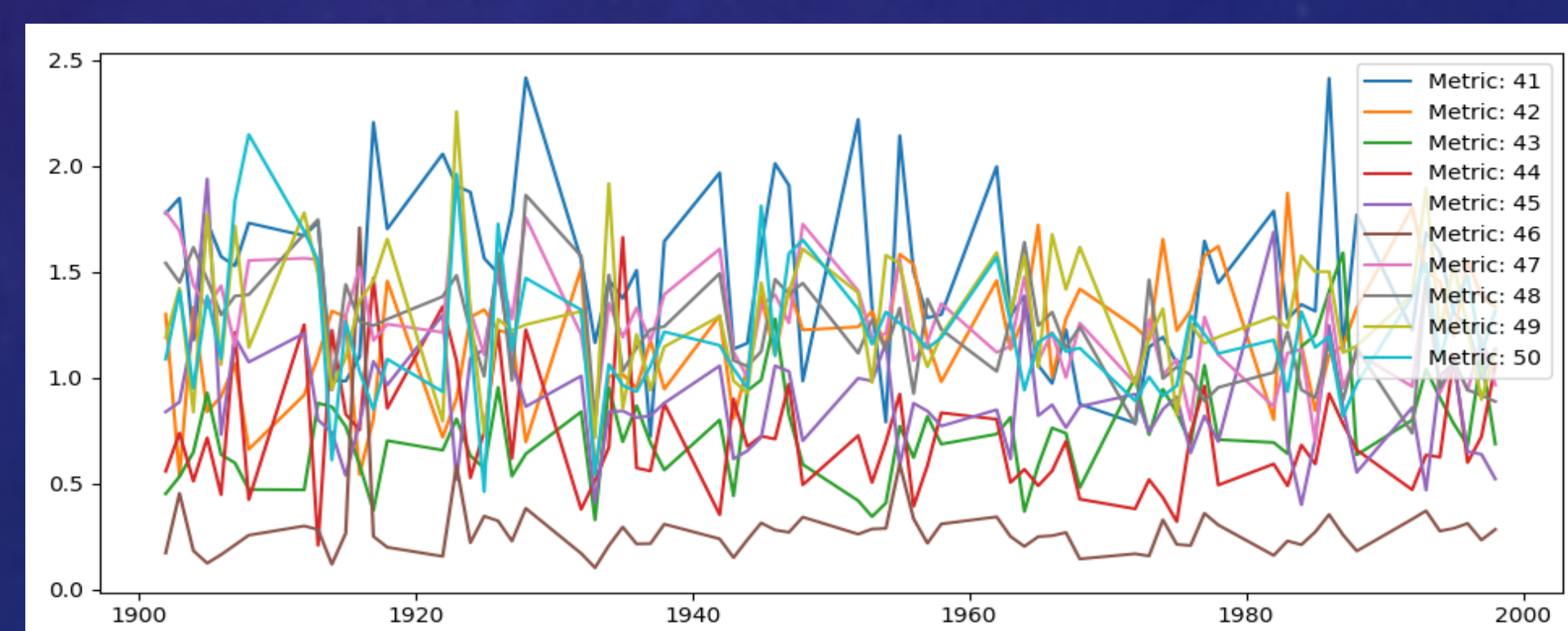
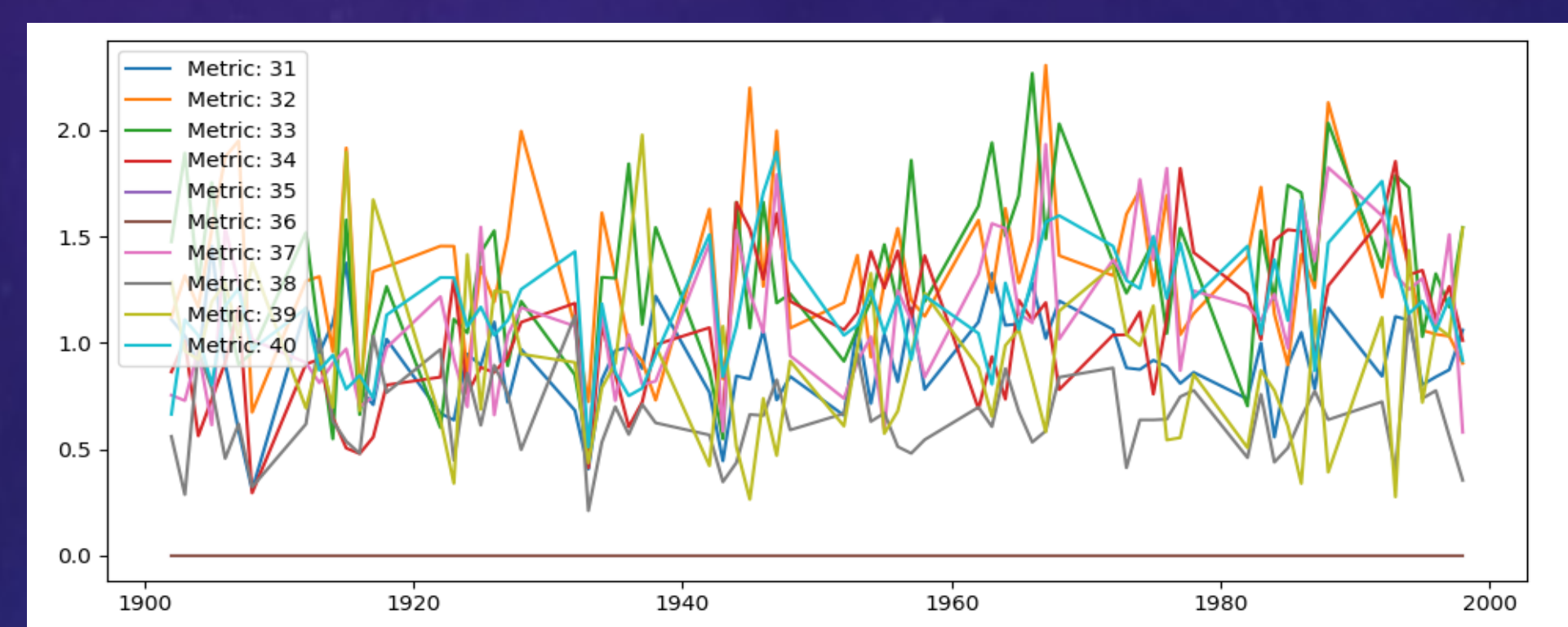
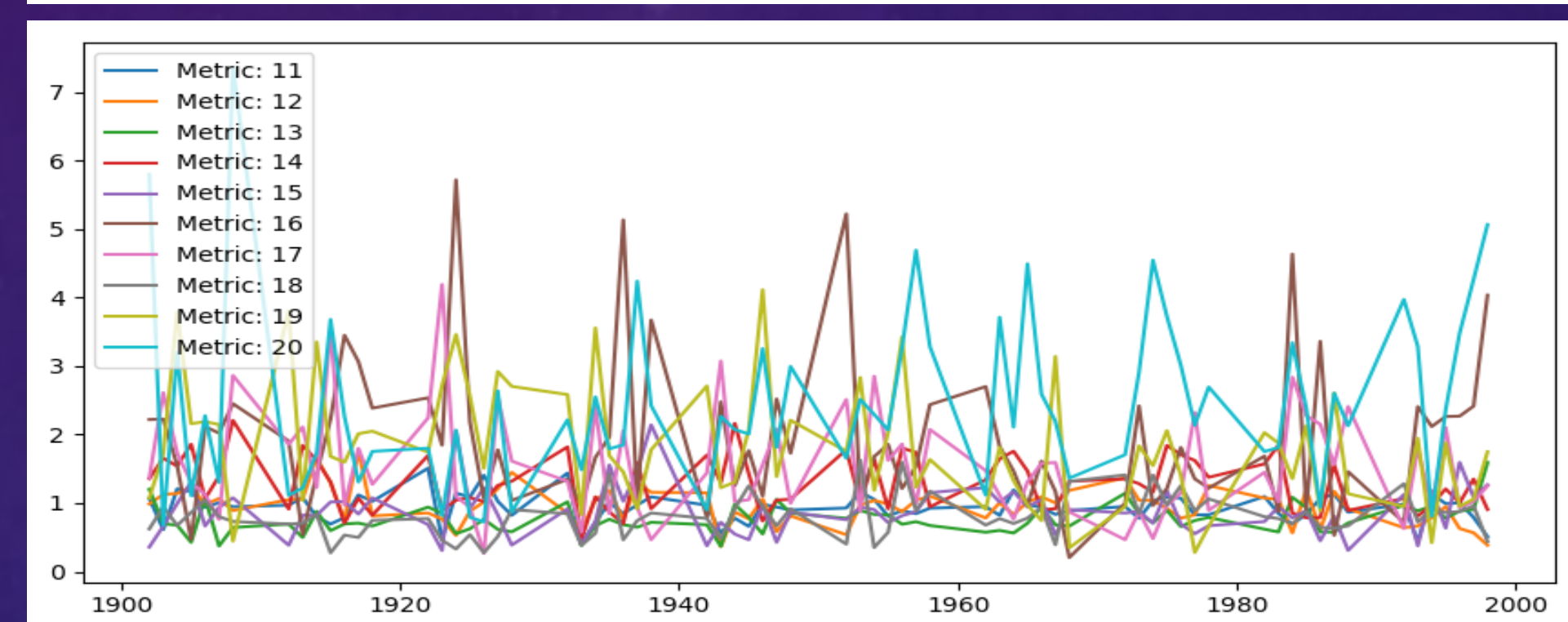
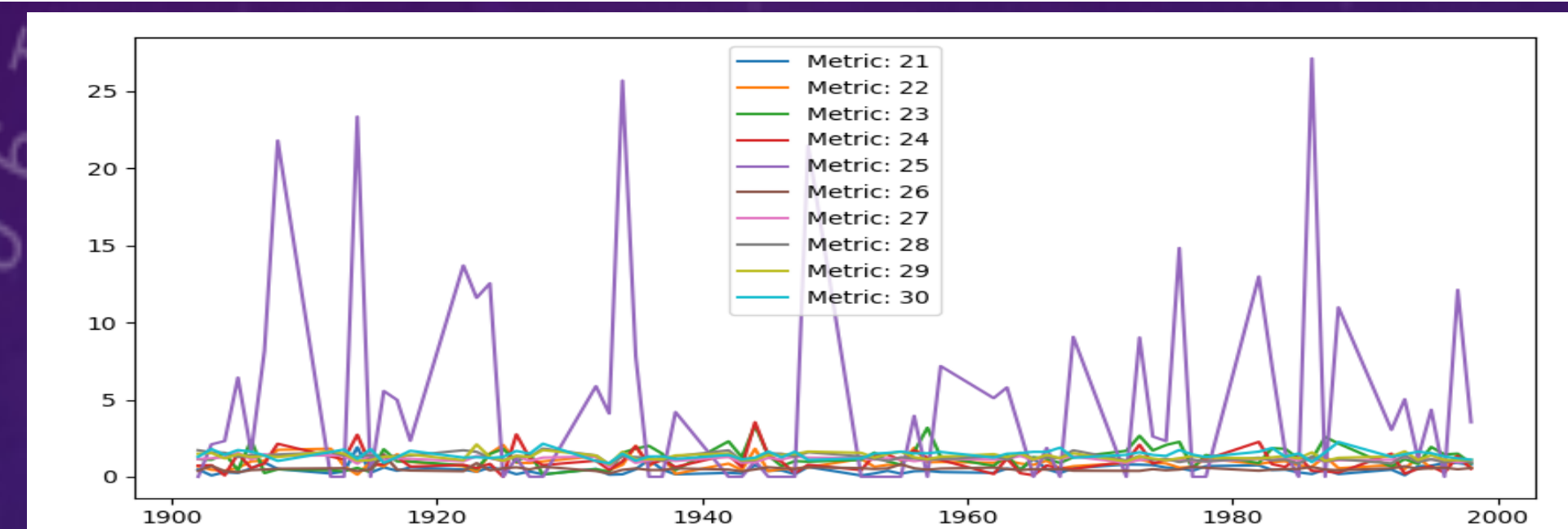
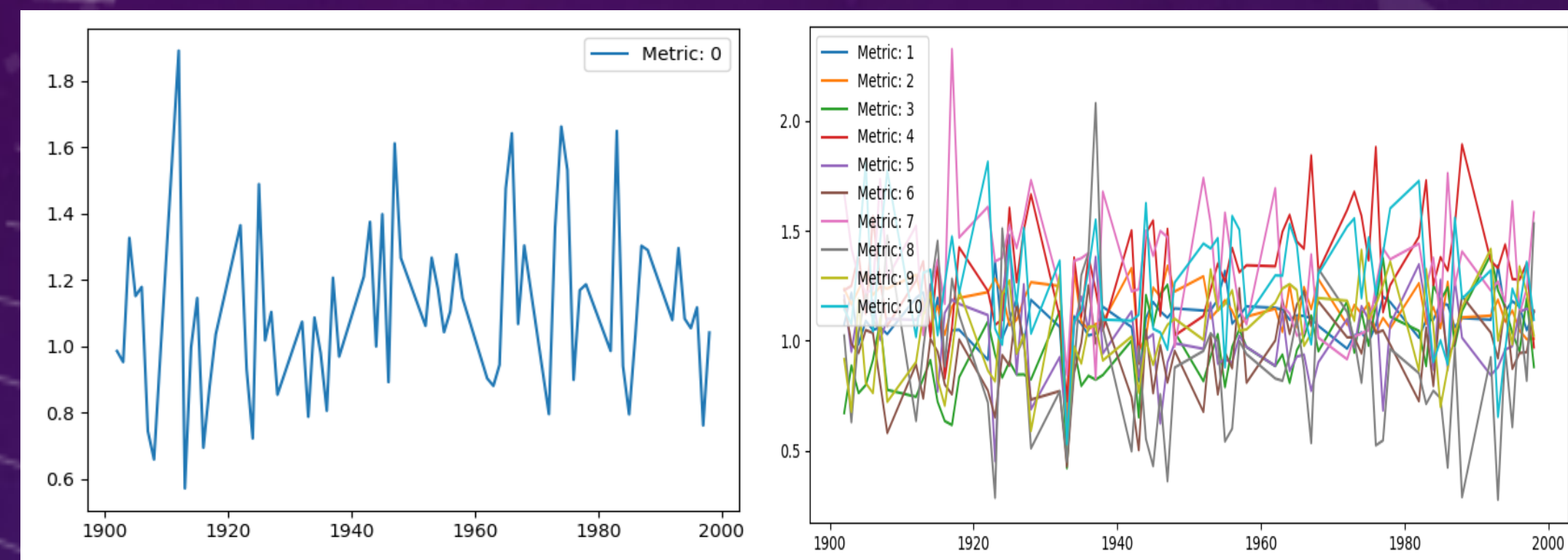
## Confusion Chart

To visualize how many source clusters were correctly associated to their proper time-span clusters, a confusion chart was made. This operates by having an x axis of actual time-span values and a y axis of predicted time-span values. Tallies are placed in each cross section showing the number of actual time spans placed into the predicted time span cluster. Ideally, high tallies should be along the diagonal of the confusion chart indicating the correct predicted time-spans were placed in the proper actual time-span slots



## Broader Impact

If successful, this research can provide insight as to how American communication has changed over time and model how it may continue to change over time as well.

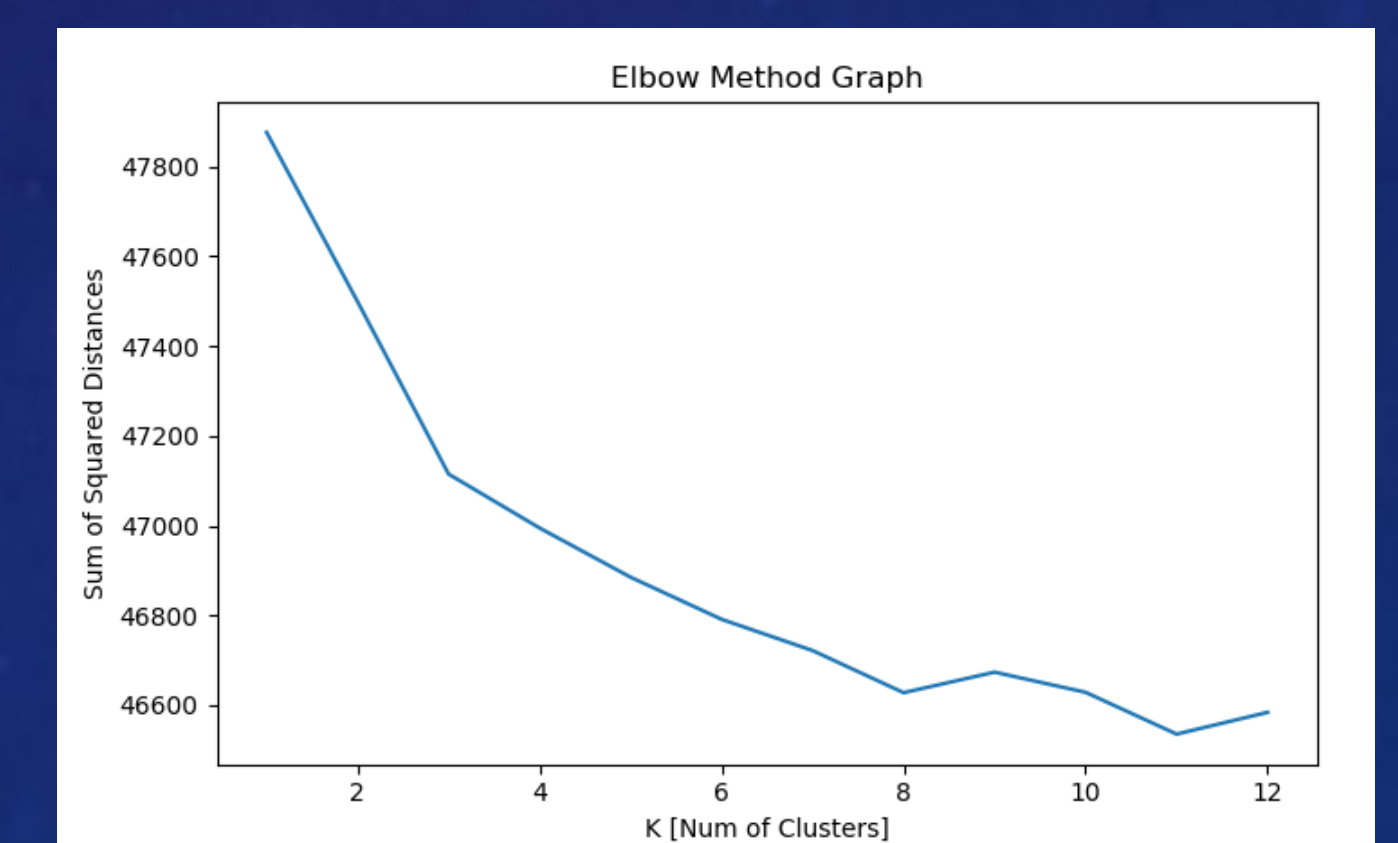


# Clustering

To accomplish prediction, data clustering is used. Clustering uses several data points, each representing a different text source. In the instance of this research the clusters are 66 dimensional points and grouped in close proximity with an association to a single cluster group. This process is highly iterative and continuously changes the location of the mother cluster points until the distance to all of its associated points is nominal. Essentially, clustering groups points with similar metrics and is used to represent a different time-span once the iterative process is completed.

## Elbow Method

To determine how many clusters, and therefore time-spans, the elbow method was used. The elbow method performs data clustering for a number of clusters and determines the sum of squared distances for each cluster and its respective points. When these distances are graphed opposed to the number of clusters used, a trend will arise where the distance decreases from a near vertical slope and levels out after a number of clusters, forming what looks to be an elbow.



## Results

Due to the limited corpora used, prediction beyond precision of an indiscriminate projection has not yet been achieved in this research project.

## Special Acknowledgments

## Dr. Claude Willan

## Director of the UH Digital Research Commons

UNIVERSITY of  
**HOUSTON**